

Wright State University

CORE Scholar

Kno.e.sis Publications

The Ohio Center of Excellence in Knowledge-
Enabled Computing (Kno.e.sis)

2007

Schema-Driven Relationship Extraction from Unstructured Text

Cartic Ramakrishnan

Wright State University - Main Campus

Follow this and additional works at: <https://corescholar.libraries.wright.edu/knoesis>



Part of the [Bioinformatics Commons](#), [Communication Technology and New Media Commons](#), [Databases and Information Systems Commons](#), [OS and Networks Commons](#), and the [Science and Technology Studies Commons](#)

Repository Citation

Ramakrishnan, C. (2007). Schema-Driven Relationship Extraction from Unstructured Text. .
<https://corescholar.libraries.wright.edu/knoesis/266>

This Presentation is brought to you for free and open access by the The Ohio Center of Excellence in Knowledge-Enabled Computing (Kno.e.sis) at CORE Scholar. It has been accepted for inclusion in Kno.e.sis Publications by an authorized administrator of CORE Scholar. For more information, please contact library-corescholar@wright.edu.

Schema-Driven Relationship Extraction from Unstructured Text

Cartic Ramakrishnan
Kno.e.sis Center, Wright State University,
Dayton, OH

Motivation

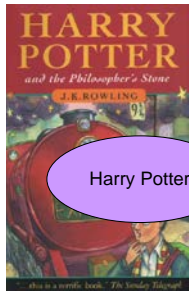
Problem Description & Approach

Results

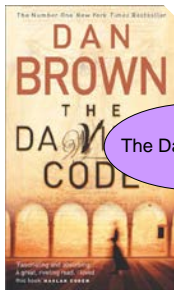
Future Work

UNDISCOVERED PUBLIC KNOWLEDGE

Discovering connections hidden in text



Harry Potter



The Da Vinci code

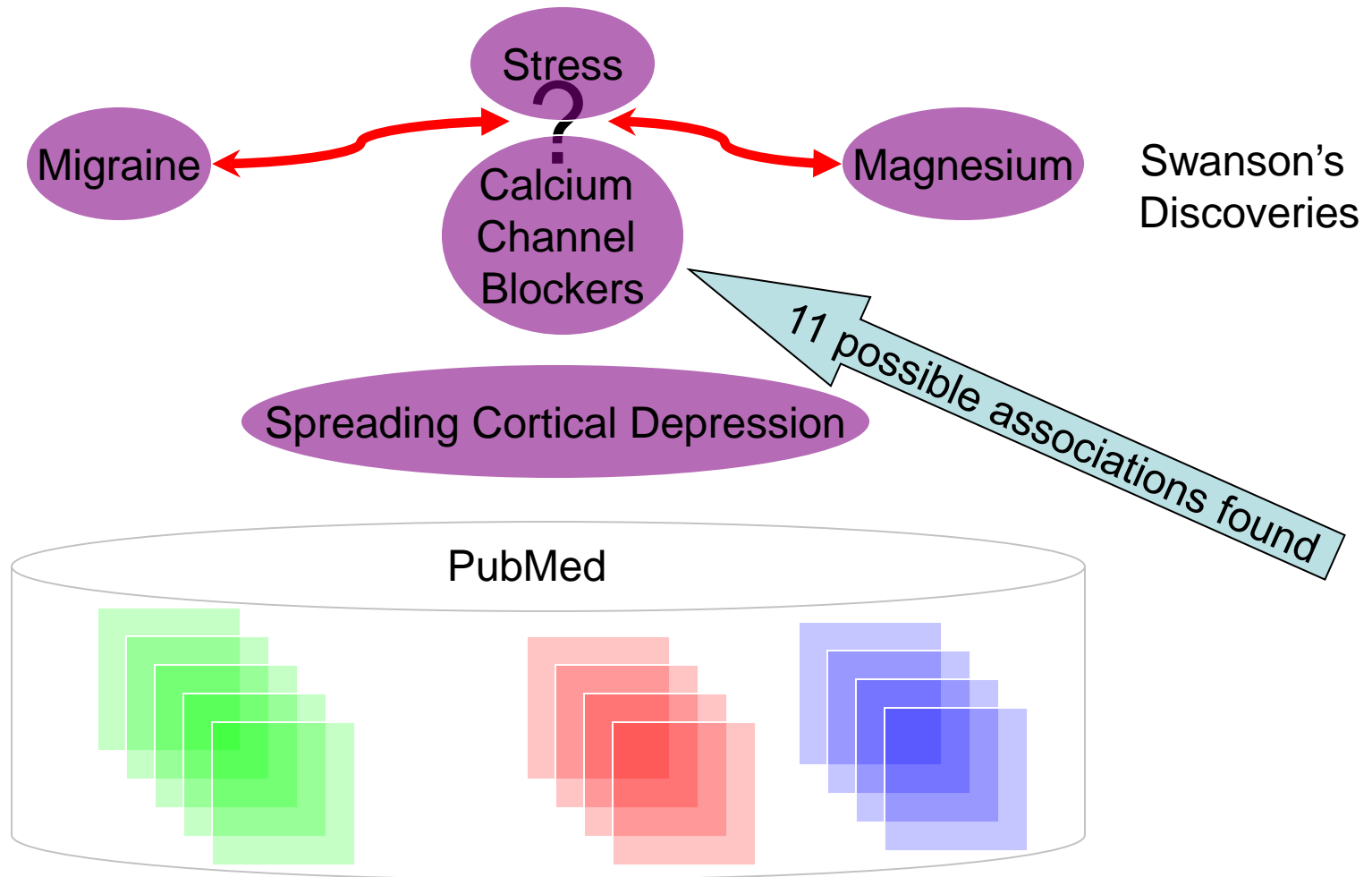


Et in Arcadia Ego



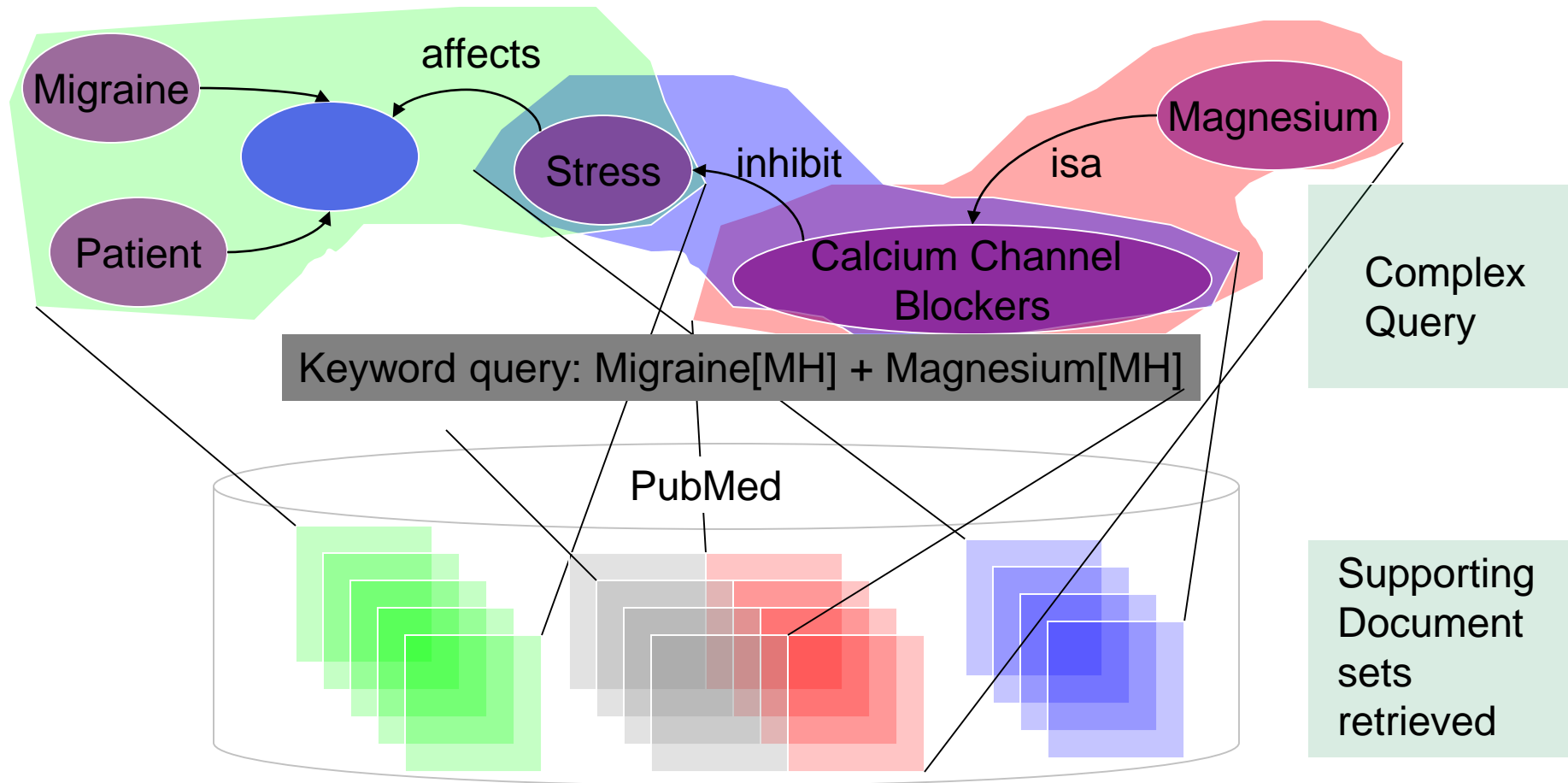
Santa Maria delle Grazie

Motivation 1 – Undiscovered Public knowledge in biology



Associations Discovered based on keyword searches
These associations were discovered in 1986
followed by manually analysis of text to establish possible relevant relationships

Motivation 2 - Hypothesis Driven retrieval of Scientific Literature



Data captured per year = 1 exabyte (10^{18})

(Eric Neumann, Science, 2005)

How much is that?

- Compare it to the estimate of the total words ever spoken by humans = 12 exabyte

A small but significant portion is text data

- PubMed 16 Million abstracts
- MedlinePlus – health information
- OMIM – catalog of human genes and genetic disorders

Undiscovered public knowledge may have also increased by a large amount

Semantic Associations over RDF graphs

– Discovery and Ranking

It is therefore critical to bridge the gap between unstructured and structured data by extracting entities and relationships between resulting in semantic metadata

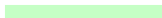
Motivation

Problem Description & Approach

Results

Future Work

Problem – Extracting relationships between MeSH terms from PubMed



UMLS – A high level schema of the biomedical domain

- 136 classes and 49 relationships
- Synonyms of all relationship – using variant lookup (tools from NLM)
- 49 relationship + their synonyms = ~**350** mostly verbs

MeSH

- 22,000+ topics organized as a forest of 16 trees
- Used to query PubMed

T147—effect
T147—induce
T147—etiology
T147—cause
T147—effecting
T147—induced

PubMed

- Over 16 million abstract
- Abstracts annotated with one or more MeSH terms

[1254239-1] An excessive endogenous or exogenous stimulation by estrogen induces adenomatous hyperplasia of the endometrium.

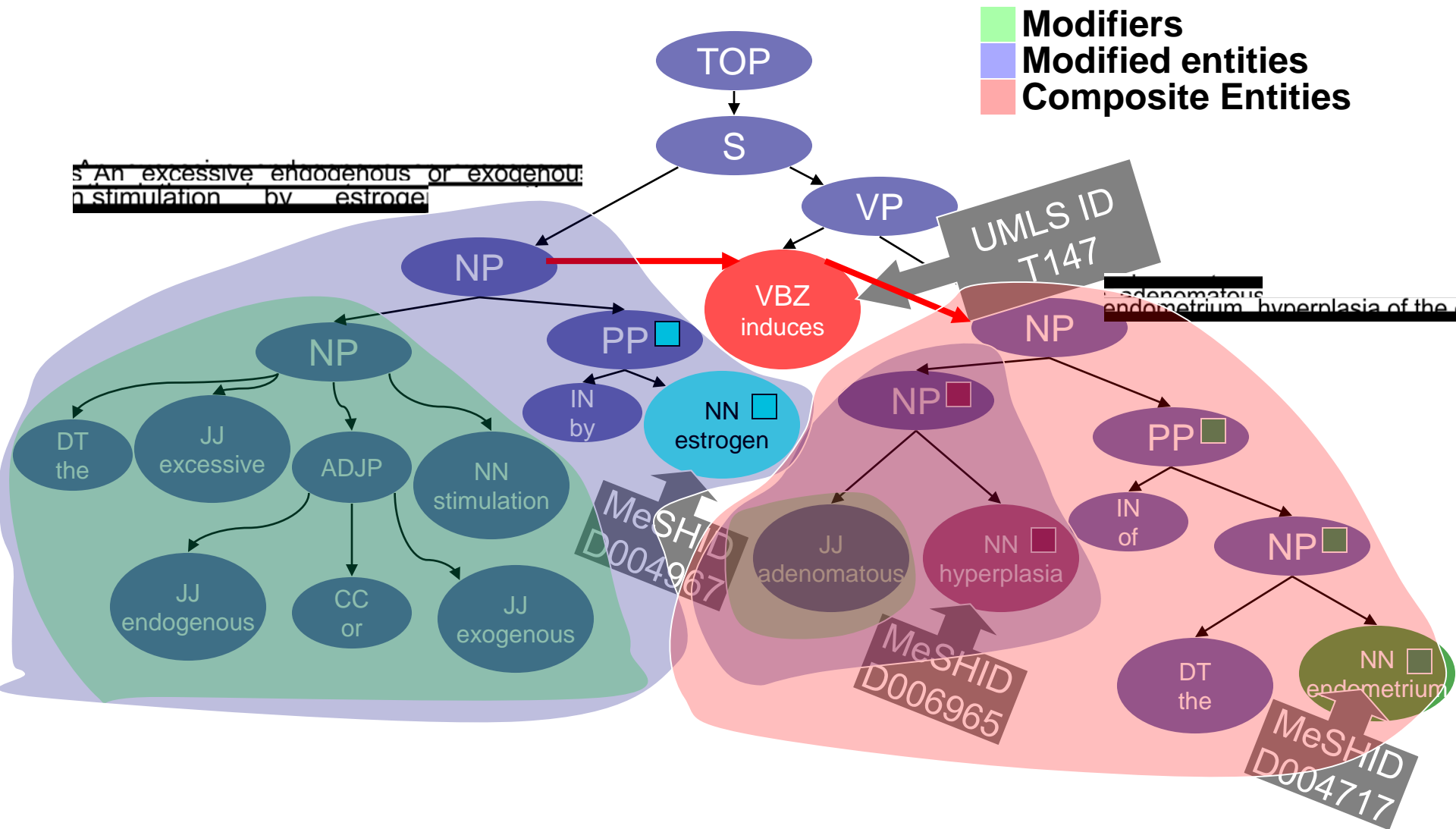


- Entities (MeSH terms) in sentences occur in modified forms

(TOP (\$adenomatous\$)NP (DT An)JJ excessivehyperplasiaADJP (JJ endogenous) (CC or) (JJ exogenous)) (NS stimulation) (PP (IN by) (NP (NN estrogen))) (VP (VBZ induces)NP (JJ adenomatous) (NN hyperplasia)) (PP (IN of) (NP (DT the) (NN endometrium))) as composites of 2 or more other entities

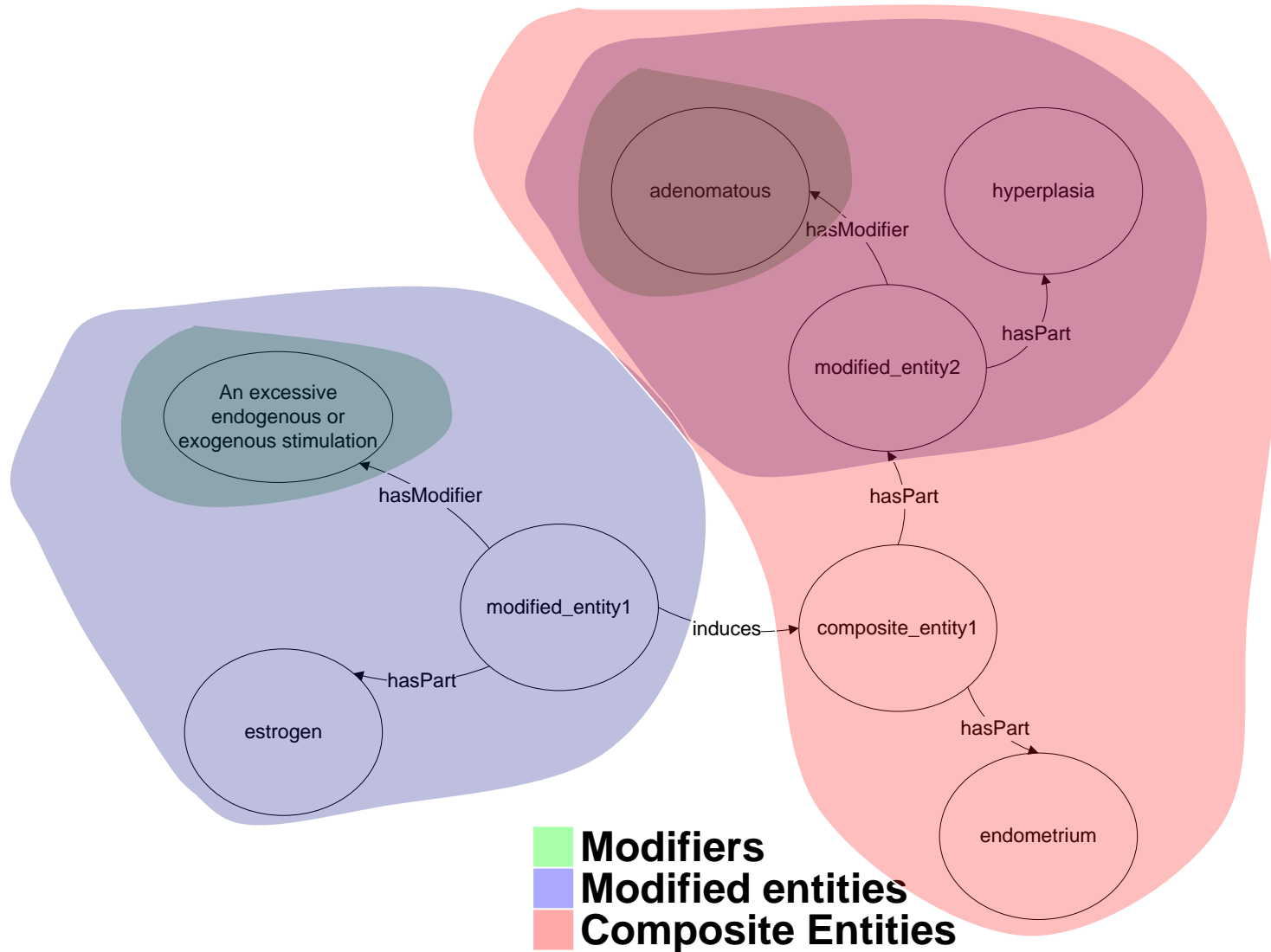
- “adenomatous hyperplasia” and “endometrium” occur as “adenomatous hyperplasia of the endometrium”

Method – Identify entities and Relationships in Parse Tree



To capture the various types of entities we define

- Simple entities as MeSH terms
- Modifiers as siblings of entities that are
 - Determiners – “Y induces **no** X”
 - Noun Phrases – “An excessive endogenous or exogenous stimulation”
 - Adjective phrases – “adenomatous”
 - Prepositional phrases – “M is induced **by the X in the Z**”
- Modified Entities as **any** entity that has a sibling which is a modifier
- Composite Entity as **any** entity that has another entity as a sibling



Motivation

Approach

Results

Future Work

Dataset 1

– Swanson's discoveries

- Associations between Migraine and Magnesium [Hearst99]
 - stress is associated with migraines
 - stress can lead to loss of magnesium
 - calcium channel blockers prevent some migraines
 - magnesium is a natural calcium channel blocker
 - spreading cortical depression (SCD) is implicated in some migraines
 - high levels of magnesium inhibit SCD
 - migraine patients have high platelet aggregability
 - magnesium can suppress platelet aggregability

Keywords pairs e.g. stress + migraine etc. against PubMed return PubMed abstracts that are annotated (by NLM) with both terms

8 pairs of terms in this scenario result in 8 subsets of PubMed

Semantic Metadata

- Represented in RDF
- With complex entities and relationships connecting them
- Pointers to original document and sentence
- Size
 - ~2MB RDF for Migraine Magnesium subset of PubMed

Ideal method to evaluate the Extraction method

- Domain experts read a set of abstract given a set of relationship names and entities to look for
- In addition to this give them the extracted triples and entities
- For every abstract the expert judges counts the correct, incorrect and missed triples
- Measure precision and recall

In the absence of a domain expert we focus of getting a feel for the utility of the extracted data

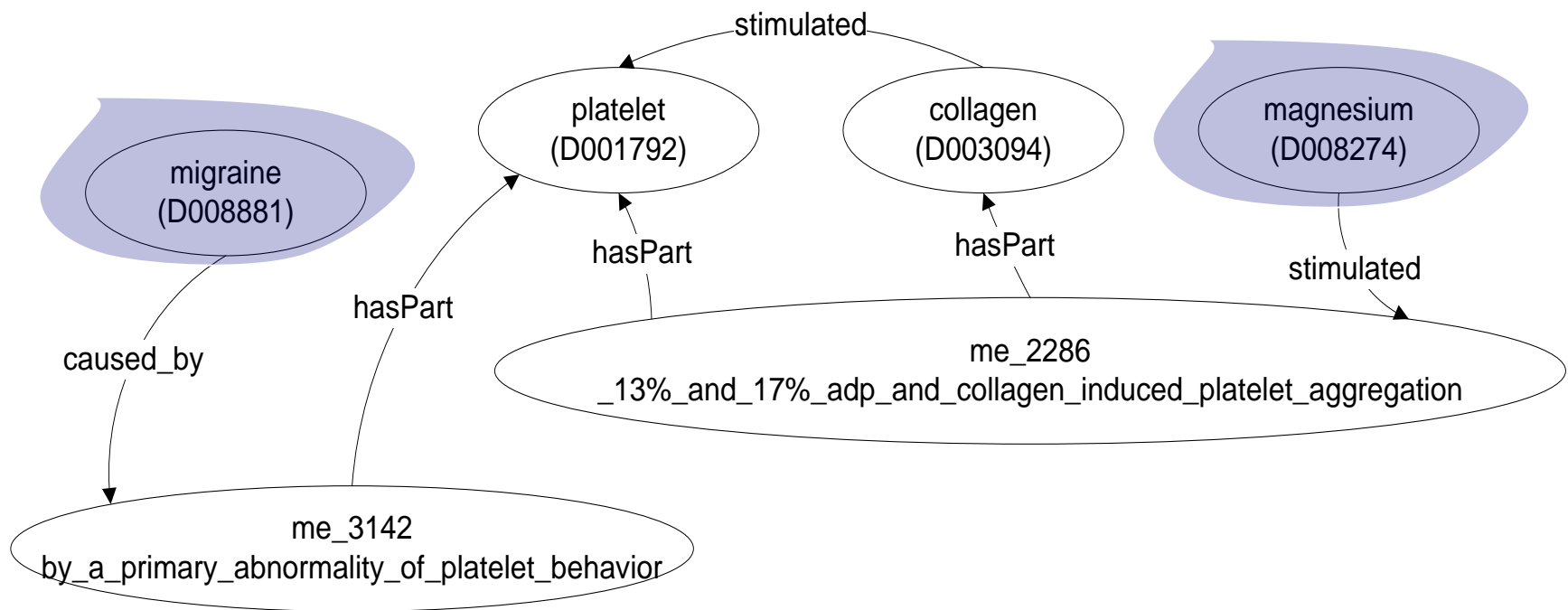
- We know the association manually discovered between Migraine and Magnesium
- We locate paths of various lengths between them and manually inspect these paths
- If the paths are indicative of the manually discovered associations the extracted data is useful

Table 2. Paths between Migraine and Magnesium

Paths between Migraine and Magnesium			
Path length	Total Number of paths found	# of interesting paths	Max. # of named relationships in any path
6	260	54	4
8	4103	1864	5
10	106450	33403	5

Paths are considered interesting if they have one or more named relationship
Other than ***hasPart*** or ***hasModifiers*** in them

An example of such a path



Dataset 2

- Neoplasm (C04)
 - For subtree of MeSH rooted at Neoplasms all topics under this subtree are used as query terms against PubMed
 - The resulting dataset contains ~500,000 PubMed abstracts
 - The extraction process run on this data returns ~150MB

Processing the tagged and parsed sentences for Dataset 2 (Neoplasm) to generate RDF took approx. **5 minutes**

Stats

- 211 different named relationships found
- 500,000 instance-property-instance statements
- 260,000 instance-property-literal statements

Currently setting up to extract RDF from all of PubMed

Motivation

Problem Description & Approach

Results

Future Work

Short-term goals (1 month)

- MeSH qualifiers (blood pressure, contraindications)
- **Curate** and release Migraine-Magnesium RDF

Long-Term goals

- More complex structures
 - Conjunctions
 - X causes Y to inhibit Z
- Rule-action language to test new extraction rules
- Finding new terms to enrich existing vocabularies
- Perhaps ontology enrichment

From ...

Hypothesis driven “wet lab” experiments

To ...

Data-driven reduction/pruning of hypothesis space leading to new insight and possibly discovery

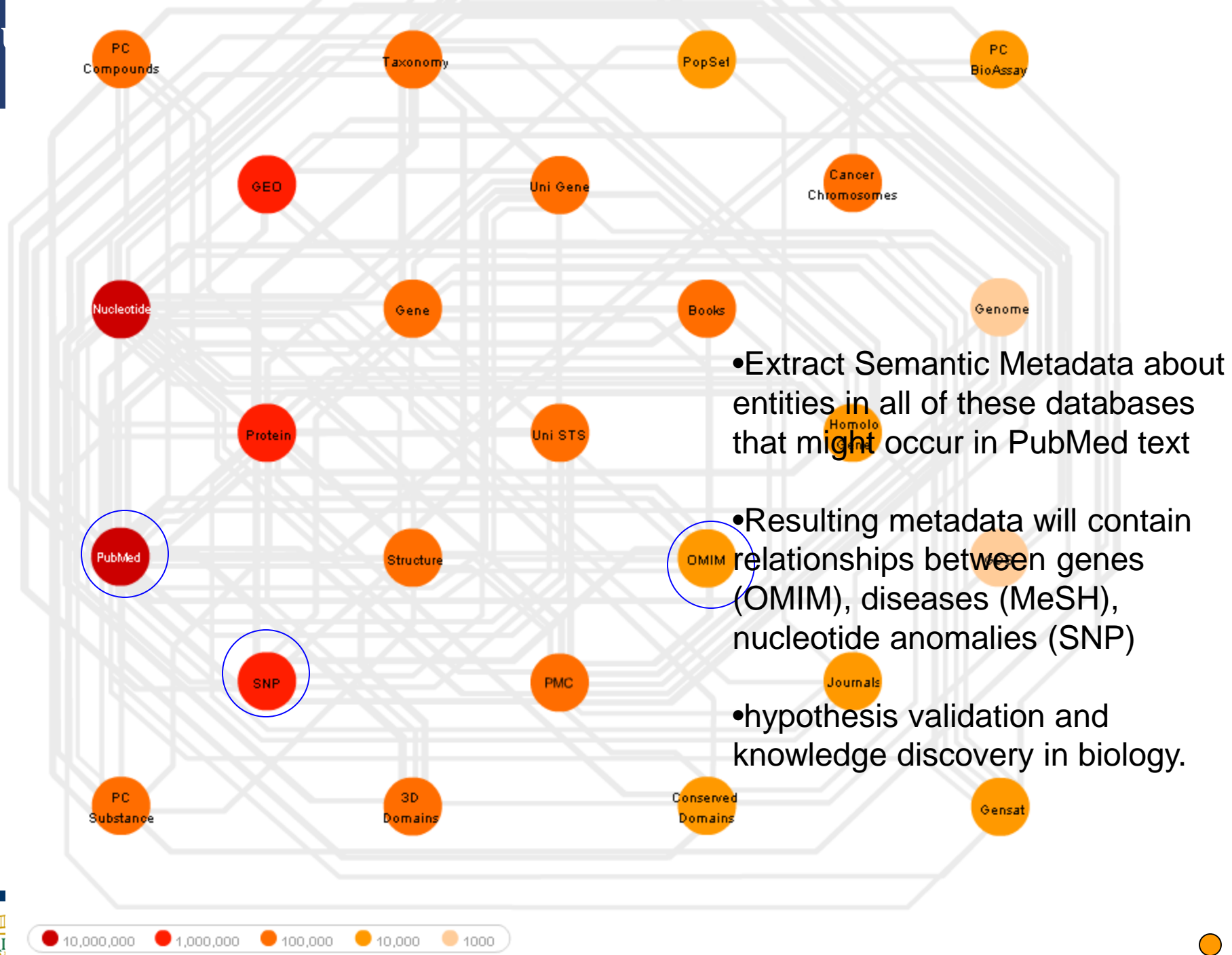
What challenges does this transition bring?

Semantic Browsing of PubMed based on named relationships between MeSH terms

Path/hypothesis based document retrieval

Knowledge discovery from literature

- Coprus-based complex relationship discovery and ranking
- Corpus-based relevant connection subgraph discovery



THANK YOU!